

Ideal gas behavior of rotamerically defined conformers in native globular proteins

Kai Wang,[†] Shiyang Long,[†] Zhiming Zhang,[†] Lanru Liu,[†] Qimeng Wang,[†] and Pu Tian^{*,†,‡}

*College of Life Science, and MOE Key Laboratory of Molecular Enzymology and Engineering
Jilin University
2699 Qianjin Street, Changchun 130012*

E-mail: tianpu@jlu.edu.cn

Abstract

Protein side chain entropy has been found to be important by both experimental and computational studies. However, the connection between side chain torsional states and protein conformational distributions remains vague. Based on the robustness of side chain rotameric states observed in both experimental structures and large scale molecular dynamics simulations, we propose to define unique combinations of side chain rotameric states as basic conformers, termed RCONFs, for entropy calculation. we further hypothesize that all RCONFs have the same constant local configurational integral for a given protein under specified solvent conditions. It follows from this hypothesis that RCONFs behave like ideal gas configurations, RCONF based conformational entropy may be effectively expressed as $S = \ln W$, with W being the number of RCONFs that are thermally accessible, and change of free energy between two

*To whom correspondence should be addressed

[†]College of Life Science

[‡]MOE Key Laboratory of Molecular Enzymology and Engineering

Jilin University

2699 Qianjin Street, Changchun 130012

given macrostates is equivalent to that of conformational entropy with a mere difference of a negative temperature factor. The validity of the ideal gas hypothesis and inferred property of change of both conformational entropy and free energy is tested in extensive molecular dynamics (MD) simulation trajectories of six globular proteins in native state. The advantage of the corresponding end point free energy method is discussed.

Introduction

While theoretical importance of entropy in physical systems has been appreciated for a long time,¹⁻³ experimental evidence for decisive roles of conformational entropy in bimolecular interactions appeared only recently.⁴⁻⁷ This is due to the fact that decomposition of free energy into contributions from system comprising components and their correlations is extremely challenging. The total free energy change of a typical protein-ligand system may be written as:

$$\Delta G = \Delta H - T\Delta S \quad (1)$$

$$\Delta H = \Delta H_{pp} + \Delta H_{ll} + \Delta H_{vv} + \Delta H_{pl} + \Delta H_{lv} + \Delta H_{pv} \quad (2)$$

$$\Delta S = \Delta S_{pp} + \Delta S_{ll} + \Delta S_{vv} + \Delta S_{cross} \quad (3)$$

With letters p , l , and v in subscript represent protein, ligand and solvent respectively. Repeated subscripts (pp , ll and vv) represent molecular interactions in the same type of molecules in enthalpic terms and corresponding entropy in entropic terms. Mixed subscripts (pv , pl and lv) represent molecular interactions between different type of molecules in enthalpic terms. ΔS_{cross} is the change of entropy reduction due to correlation of molecular degrees of freedom (DOFs) between (among) different types of molecules in entropic terms^a. Computationally, systematic and quantitative evaluation of numerous methods for direct calculation of macromolecular entropy and solvent entropy

^a S_{pp} is the entropy of protein molecules under the free energy landscape determined by both the intra- and inter-molecular interactions existing under given solvent conditions. S_{ll} and S_{vv} are defined similarly. When ligands are unbound, $S_{pl} \approx 0$ in dilute solution and $S_{cross} = S_{pv} + S_{lv}$. When ligands are bound, there is actually a third order correlation term S_{plv} so $S_{cross} = S_{pl} + S_{pv} + S_{lv} + S_{plv}$.

is yet to be done.⁸⁻¹³ Furthermore, no general method is available for direct computation of correlation entropies between solute and solvent molecules (e.g. protein and water). Experimentally, the well-established isothermal calorimetry (ITC)¹⁴ measurement only directly gives ΔH and ΔG . The “model free approach” of Lipari and Szabo^{15,16} provides a theoretical connection between local motion of bond vectors, which is measurable by NMR relaxation and characterizable by the square of the generalized order parameter, and thermodynamics. In a series of NMR studies,^{6,7} change of protein configurational (conformational) entropy as manifested by side chain motion of methyl bearing residues was found to exhibit strong linear correlation with total entropy change upon binding of ligands. Consequently, the ensemble average of side chain entropies of methyl bearing residues may effectively serve as an “entropy meter” for binding entropy. This strategy was successfully utilized by Tzeng *et. al.*⁵ in elegantly designed CAP (catabolite activator protein)-DNA interaction systems to demonstrate the importance of configurational (conformational) entropy contribution to the free energy. Despite these significant advances, the connection between protein side chain rotameric states and thermodynamics of protein conformational distributions remains to be constructed.

The ideal gas hypothesis of rotamerically defined conformers

For a protein molecule, a binding event with a ligand is essentially a redistribution of equilibrium statistical weight in the conformational space. The observed linear correlation between entropy change of methyl bearing side chains and that of ligand binding suggests that there might be undiscovered general rule(s) governing distributions of rotamerically defined conformers (RCONFs) in the conformational space. To be more specific on the definition of RCONFs, each RCONF is defined by a given combination of torsional states of unique side chain all-heavy-atom torsional DOFs. Two structural states of a protein belong to the same RCONF if and only if they share the same torsional state for each unique side chain all-heavy-atom torsional DOF.

The configurational integral of a protein molecule in solution may be written as:

$$Z = \int e^{-\beta U(r_p, r_v)} dr_p dr_v \quad (4)$$

with β being the reciprocal temperature, U being the potential energy, and r_p and r_v being coordinates of protein and solvent atoms. If we partition the whole conformational space into small basic elements, each of which corresponds to a unique RCONF, the configurational integral can be written as a sum of all sub-integrals (local integrals) corresponding to RCONFs:

$$Z = \sum_{i=1}^{N_{rconf}(U(r_p, r_v))} Z_i \quad (5)$$

$$Z_i = \int_{rconf_i} e^{-\beta U(r_p, r_v)} dr_p dr_v \quad (6)$$

For two arbitrarily given macrostates A and B :

$$Z^A = \sum_{i=1}^{N_{rconf}^A(U(r_p, r_v))} Z_i^A \quad (7)$$

$$Z^B = \sum_{j=1}^{N_{rconf}^B(U(r_p, r_v))} Z_j^B \quad (8)$$

$$Z_i^A = \int_{rconf_i^A} e^{-\beta U(r_p, r_v)} dr_p dr_v \quad (9)$$

$$Z_j^B = \int_{rconf_j^B} e^{-\beta U(r_p, r_v)} dr_p dr_v \quad (10)$$

$$\Delta F^{AB} = k_B T \ln \frac{Z^A}{Z^B} \quad (11)$$

ΔF^{AB} is the change of Helmholtz free energy between macrostates A and B . If we consider limit of dilute solution and neglect change of rotational and translational entropy, then ΔS_{pp} in equation (3) is reduced to change of configurational entropy ΔS_{config} . Based on the definition of RCONFs and following a common practice in theoretical studies,¹⁷⁻²⁰ we may split configurational entropy S_{config} into conformational and vibrational contributions (a physics based proof is provide

by Chang and Gilson¹⁹) as shown below:

$$S_{config} = S_{rconf} + S_{rconf-vib} \quad (12)$$

$$S_{rconf} = -k_B \sum_{i=1}^{i=N_{rconf}(U(r_p, r_v))} P_i \ln P_i \quad (13)$$

$$S_{rconf-vib} = \sum_{i=1}^{i=N_{rconf}(U(r_p, r_v))} P_i S_{rconf-vib}^i \quad (14)$$

For two arbitrarily given macrostates (conformation) A and B ,

$$S_{config}^A = S_{rconf}^A + S_{rconf-vib}^A \quad (15)$$

$$S_{rconf}^A = -k_B \sum_{i=1}^{i=N_{rconf}^A(U(r_p, r_v))} P_i^A \ln P_i^A \quad (16)$$

$$S_{rconf-vib}^A = \sum_{i=1}^{i=N_{rconf}^A(U(r_p, r_v))} P_i^A S_{rconf-vib}^{A_i} \quad (17)$$

$$S_{config}^B = S_{rconf}^B + S_{rconf-vib}^B \quad (18)$$

$$S_{rconf}^B = -k_B \sum_{j=1}^{j=N_{rconf}^B(U(r_p, r_v))} P_j^B \ln P_j^B \quad (19)$$

$$S_{rconf-vib}^B = \sum_{j=1}^{j=N_{rconf}^B(U(r_p, r_v))} P_j^B S_{rconf-vib}^{B_j} \quad (20)$$

$$\Delta S_{config}^{AB} = \Delta S_{rconf}^{AB} + \Delta S_{rconf-vib}^{AB} \quad (21)$$

S_{rconf} is the conformational entropy based on our definition of RCONFs, $S_{rconf-vib}$ include contributions from both bonding and bending vibrational entropies, and from local part of torsional entropy within each RCONF, and $P_{i(j)}^{A(B)}$ is the probability of the $i(j)$ th RCONF in the macrostate $A(B)$. The rational of only considering side chain torsional DOFs will be discussed after presentation of the main results. As is apparent from equations (12-14), allocation of S_{config} between the “conformational term” and the “vibrational term” depends on specific definition of conforma-

tions. The robustness of rotameric states as observed in both experimental structures^{21,22} and MD simulations²³ provides a feasible and practical base for general applicability of RCONF-based conformational entropy. Before proceeding further, we need to clarify that two terms “configurational entropy” and “conformational entropy” are used interchangeably in many experimental studies.^{5–7} In this study, the term S_{rconf} denotes the entropy based on distributions of RCONFs, it does not consider details of microstate distributions within any given RCONF, and is apparently different from S_{config} , which includes all vibrational contributions.

For typically-sized natural proteins, N_{rconf} is an astronomically large number. For example, with fixed backbone and rigid rotamer model, the allowed combination of side chain rotameric states is estimated to be as many as $\sim 10^{40}$ for 76-residue ubiquitin given the backbone coordinates in the PDB code *1ubq*. Allowing both side chain and backbone flexibility is likely to increase N_{rconf} , which probably increase exponentially with the number of residues and consequently will certainly be more intractable for larger proteins. Therefore, solving local integrals Z_i s within RCONFs is a not a feasible path for calculating configurational integral Z of the whole conformational space, or Z^A of a given macrostate A . A drastic simplification is to assume that all Z_i s for a given protein under specific solvent conditions have the same constant value Z_{rconf} across the whole conformational space, and consequently:

$$Z = Z_{rconf} N_{rconf}(U(r_p, r_v)) \quad (22)$$

$$\Delta F^{AB} = k_B T \ln \frac{N_{rconf}^A(U(r_p, r_v))}{N_{rconf}^B(U(r_p, r_v))} \quad (23)$$

$$S_{rconf} = k_B \ln N_{rconf}(U(r_p, r_v)) \quad (24)$$

$$\Delta S_{rconf}^{AB} = k_B \ln \frac{N_{rconf}^B(U(r_p, r_v))}{N_{rconf}^A(U(r_p, r_v))} \quad (25)$$

Under this drastically simplified assumption of constant local configurational integral across all RCONFs, RCONFs of proteins behave like configurations of ideal gas, and change of free en-

ergy between two macrostates is equivalent to that of S_{rconf} with a mere difference of a negative temperature factor. It is important to note that the hypothesized constant Z_{rconf} is dependent on both identity of protein molecules and specific solvent conditions. Additionally, even under this assumption, $S_{rconf-vib}$ for different RCONFs may vary. Since it is demonstrated that entropic contributions of hard DOFs (bonding and bending) is separable from that of soft (torsional) ones,²⁴ and our primary goal is to investigate how RCONFs distribute in the protein conformational space, we leave out the complexity of analyzing bonding and bending vibrational entropies in this study.

The absolute RCONF-based conformational entropy S_{rconf} of a typically-sized protein is extremely difficult, if ever possible, to be obtained by a direct sampling approach, such as MD simulations or *regular biochemical experiments*, where usually micro-molar or much less proteins are used. For the above mentioned example of ubiquitin with backbone configurations fixed as the in the crystal structure *1ubq*, assuming a 10-*fs* average life time and that each RCONF is visited only once, it will take 1 mole of ubiquitin (~ 9 kilograms) a few minutes to complete a traverse of $\sim 10^{40}$ RCONFs. In reality, the measured change of conformational entropy for a typical major protein conformational change between two macrostates *A* and *B* is the change (presumably converged) of observed local conformational entropies.

$$\Delta S_{rconf}^{AB} = \Delta S_{rconf}^{O-AB} = S_{rconf}^{O-A} - S_{rconf}^{O-B} \quad (26)$$

with *O*— indicate observed part of the specified conformational space defining a macrostate. Again, it is important to emphasize this is not only the case for MD simulations, but also true for typical biochemical measurements and physiological activity of proteins. Since direct analysis of RCONF distributions in the whole conformational space is intractable, we take a step back and analyze S_{rconf}^{O-} for arbitrarily selected region of the visited protein conformational space. The logic is that *rule of RCONF distribution manifested among arbitrarily selected parts of conformational space should be effectively true for the whole conformational space*. MD trajectories provide a convenient path for arbitrary partitioning of conformational space for given protein molecules.

It is well established in the informational theory field²⁵ that for a static distribution with well-defined basic states, as in the case of equilibrium distributions of RCONFs in the conformational space, entropy may be constructed by arbitrary division of the whole system into M subparts.

$$S = - \sum_{i=1}^{i=N} P_i \ln P_i = - \sum_{j=1}^{j=M} P_j \ln P_j + \sum_{j=1}^{j=M} P_j S_j \quad (27)$$

$$S_j = - \sum_{k=1}^{k=k_j} P_k \ln P_k \quad (j = 1, 2, \dots, M) \quad (28)$$

$$N = \sum_{j=1}^{j=M} k_j \quad (29)$$

with P_i , P_j and P_k being properly normalized:

$$\sum_{i=1}^{i=N} P_i = 1, \quad \sum_{j=1}^{j=M} P_j = 1 \quad \text{and} \quad \sum_{k=1}^{k=k_j} P_k = 1 \quad (j = 1, 2, \dots, M) \quad (30)$$

S is the global informational entropy and S_j ($j = 1, 2, \dots, M$) are local informational entropies, it is noted that such division may be carried out recursively. We may similarly divide the whole conformational space of a protein into M arbitrary parts according to our need. In reality, we rarely care the global conformational entropy of a protein molecule. Instead, what we are most interested in are differences between local conformational entropies of relevant macrostates (conformations).

We carried RCONF analysis based on extensive MD trajectories of six globular proteins in native ensemble to test the ideal gas hypothesis. While this approximation does not perform well for absolute value of S_{rconf} for given macrostates at very fine time resolution (10fs) of observation, it is demonstrated to be highly accurate and reliable, when ΔS_{rconf} or ΔF is the major concern, for all investigated protein molecules regardless of the time resolution of observation. The advantage of the end-point free energy estimation strategy as indicated by equation (23) is discussed.

Results

Ideal gas behavior of RCONFs

To test our hypothesis, we collected extensive MD trajectories of six globular proteins with different folds and sizes. Structures of these proteins are presented in Fig. 1. By encoding torsional states of side chains into bit vectors and using the radix sorting algorithm,²⁶ we assigned snapshots from MD trajectories to unique RCONFs, which are defined according to rotameric states as reported by Scouras and Daggett,²³ the relevant information of MD trajectories and the results are listed in Table 1. For three trajectory sets BPTI-a, CDK2 and BamC, each snapshot corresponds to a unique RCONF (i.e. $N_{rconf} = N_{snap}$). Consequently, all sampled RCONFs have the same observed probability under the given time resolutions (250ps, 2ps and 2ps) for the three proteins. Therefore, the observed RCONFs in these trajectories behave like ideal gas configurations. When the visited conformational space is divided into M arbitrary partitions:

$$S_{rconf}^{O_{j-}} = -k_B \sum_{k=1}^{k=n_{rconf}^j} P_k \ln P_k = k_B \ln(n_{rconf}^j) \quad (n_{rconf}^1 + n_{rconf}^2 + \dots + n_{rconf}^M = N_{snap}) \quad (31)$$

with $S_{rconf}^{O_{j-}}$ being the observed local RCONF-based conformational entropy for the j th of the given M arbitrary partitions. n_{rconf}^j being the number of observed RCONFs in the corresponding region of conformational space. However, for other trajectory sets of different proteins (HEWL and BamE), it is apparent that N_{rconf} does not equal to N_{snap} anymore. To quantitatively characterize deviations from ideal gas behavior of recorded RCONFs in these trajectories, we calculated for each set of protein trajectories both the observed S_{rconf}^{O-} and its ideal gas approximations S_{rconf}^{O-ig} for the whole visited conformational space as shown in the equation below:

$$S_{rconf}^{O-} = -k_B \sum_{i=1}^{i=N_{rconf}} P_i \ln P_i \quad (32)$$

$$S_{rconf}^{O-ig} = k_B \ln(N_{rconf}) \quad (33)$$

the deviation from the ideal gas approximation:

$$\delta S_{rconf} = S_{rconf}^{O-} - S_{rconf}^{O-ig} \quad (34)$$

is listed in Table 1. Non-zero δS were observed for some MD trajectories with pico-second(s) ($1ps$ to $4ps$) snapshot intervals. This observation indicates that much larger deviations from ideal gas behavior may have been observed if significantly finer intervals were utilized to record MD trajectories.

To resolve this potential concern, we generated three sets of fine resolution trajectories for proteins HEWL, BPTI and KLKA, and denoted them as HEWL-b, BPTI-c and KLKA-b respectively. Origins of these trajectories are uniformly distributed in the corresponding set of trajectories recorded with pico-second(s) ($1ps$ to $4ps$) intervals. The interval for saving snapshots is set to $10fs$, which is comparable with typical bonding vibrational cycles and is expected to capture interesting torsional transitions within the visited conformational space. At this time scale resolution, N_{rconf} is indeed quite different from N_{snap} and significantly larger δS s are observed (Table 1.). We further examined the statistical weight (w_{rconf}) of RCONFs, which under the assumption of the constant local integral Z_{rconf} should be a constant across all RCONFs. As shown in Fig. 2, probability of RCONFs for HEWL follows approximately an exponential decay as a function of w_{rconf} . These observations suggest that the ideal gas hypothesis might not be helpful in dealing with protein conformational distributions.

However, what we care the most is the change of observed local conformational entropy (ΔS_{rconf}^{O-}) and free energy between macrostates in cases of interested events (e.g. conformational change or molecular binding). To analyze behavior of ΔS_{rconf}^{O-} between arbitrary partitions of conformational space visited by $10-fs$ interval trajectories of HEWL, we chose the following different ways of conformational space division. Firstly, we take a given backbone dihedral (ϕ or ψ) as the order parameter and divide the whole visited conformational space into 20 windows on it. S_{rconf}^{Oj-} and S_{rconf}^{Oj-ig} ($j = 1, 2, \dots, 20$) were calculated for each window. For each backbone torsion, such

division and calculation was performed, and a $S_{rconf}^{O_j^-}$ vs. $S_{rconf}^{O_j^{-ig}}$ plot was generated. Strong linear correlations were observed for all 256 plots. After performing linear fit, the distributions of slopes and correlation coefficients are shown in Fig. 3. Essentially, both the slope and correlation coefficient are approximately equal to 1, indicating the robustness of the ideal gas behavior as far as the change of S_{rconf} is concerned.

The above mentioned plots are constructed for sets of non-overlapping and complete conformational partitions. A set of n partitions ($\Omega_i, i = 1, 2, \dots, n$) in a specified configurational space Ω are non-overlapping and complete if:

$$\Omega_i \cap \Omega_j = \emptyset \quad (i \neq j \quad \text{and} \quad i, j = 1, 2, \dots, n) \quad (35)$$

$$\text{and} \quad \Omega = \Omega_1 \cup \Omega_2 \cup \dots \cup \Omega_n \quad (36)$$

Since what we want to analyze are distributions of RCONFs in arbitrary conformational partitions, which certainly may overlap. In reality, it is not unusual for two interested macrostates to overlap in conformational space (e.g. ligand bound and ligand free proteins for a given protein-ligand system), and such overlap is the theoretical foundation of the well-acknowledged conformational selection mechanism.²⁷ We therefore constructed a $S_{rconf}^{O_j^-}$ vs. $S_{rconf}^{O_j^{-ig}}$ plot for all partitions based on various backbone torsions, as shown in Fig. 4. Again, both the slope and the correlation coefficient of a linear fit are approximately being 1.0.

Next we considered two alternative order parameters for conformational space partition that may not be constructed as linear combinations of backbone torsional DOFs, radius of gyration R_g and number of native contacts N_{nc} (see ref²⁸ for specific definition of native contacts of HEWL). After these two quantities were calculated for each snapshot in the 10-*fs* interval HEWL trajectories, 20 equal-width windows were created for both R_g and N_{nc} , $S_{rconf}^{O_j^-}$ vs. $S_{rconf}^{O_j^{-ig}}$ plot was constructed for the 40 conformational partitions as shown in Fig. 4. Similarly, both the slope and the correlation coefficient of a linear fit are approximately being 1.0. The observation that the slope being approximately 1.0 in these plots implies that $\Delta S^{O^-} \approx \Delta S^{O^{-ig}}$. Therefore, although ideal gas

behavior of S_{rconf} is only true for sufficiently coarse time resolution of observation, it is a very good approximation for ΔS_{rconf} regardless of time resolution of observation.

Equivalence between RCONF based conformational entropy and free energy

As shown in equations 23 and 25, equivalence between ΔS_{rconf}^{AB} and ΔF^{AB} (except the negative temperature factor) is another major conclusion of the ideal gas hypothesis. In an equilibrium canonical system, free energy difference between two given part of conformational space may be effectively calculated based upon observed populations.

$$\Delta F^{AB} = k_B T \ln \frac{N^A}{N^B} \quad (37)$$

with A and B stands for two arbitrarily given partitions of conformational space (macrostates), and N^A and N^B being observed populations, which are effectively represented by the number of snapshots (N_{snap}^A and N_{snap}^B) in equilibrium MD trajectories. Imagining that there is a reference conformation in equilibrium with other visited parts of conformational space, and this conformation has a statistical weight corresponding to 1 snapshot, then relative free energy of any given part of conformational space A becomes $-k_B T \ln N_{snap}^A$, with N_{snap}^A being the number of snapshots in A .

We plotted relative free energy $-k_B T \ln N_{snap}$ as a function of local ideal gas entropy $S_{rconf}^{O-ig} = k_B \ln N_{rconf}$ for the above mentioned sets of conformational space partitions (see Fig. 5.) A strong linear correlation is observed with a slope of approximately -1.0 regardless of different ways of conformational space partitioning. These observations further validated the ideal gas hypotheses. To utilize equation (23) as a new end-point free energy calculation method, there are two possible paths. The first one is to directly calculate the ratio $\frac{N_{rconf}^A(U(r_p, r_v))}{N_{rconf}^B(U(r_p, r_v))}$ without knowing the absolute values of both the numerator and the denominator. This is convenient when converged equilibrium sampling of conformations A and B is readily achievable. Direct calculation of both $N_{rconf}^A(U(r_p, r_v))$ and $N_{rconf}^B(U(r_p, r_v))$ provides a potential alternative when converged sampling of both end states are difficult. However, as briefly mentioned in the introduction, due to as-

tronomically large number of RCONFs for typically sized proteins, enumeration of RCONFs is not realistic for a specified region of conformational space from regular MD simulations, from enhanced sampling techniques, or from typical biochemical experiments. On the other hand, importance sampling²⁹ in combination with sequential Monte Carlo³⁰ was demonstrated to be an efficient and reliable way of counting the number of conformers for fixed backbone and rigid side chain rotamers. With incorporation of side chain and backbone flexibility, the number of RCONFs may be effectively counted by such importance sampling procedures. The most appealing feature of this approach is that no overlapping of conformational space is required. Additionally, such calculation will be compatible with solvation treatment ranging from full explicit solvent to simple implicit solvation models and anything in between.

Similar validation analyses of ideal gas hypothesis are performed for all trajectories with a non-zero δS_{rconf} and the same conclusion is reached (data not shown).

Analysis with alternatively defined basic conformers

Most trajectories utilized in this study are generated with CHARMM22 force fields, except for the fact that CDK2 trajectories are generated with AMBER ff12SB and BPTI-a trajectories with a modified AMBER ff99SB.³¹ Differences in distributions of side chain torsional DOFs are expected between MD trajectories generated by different force fields, and between experimental structures and MD trajectories of given force fields. Based on distributions of heavy-atom-defined side chain torsional angles observed in each set of trajectories, we utilized an in-house torsional state assignment procedure²⁸ to define trajectory-set specific basic conformers, termed RCONF2s. We compared the results from HEWL trajectory set with what reported by Scouras and Daggett,²³ which was used to define RCONFs. While most boundaries for torsional state assignment only differ for a few degrees, a number of torsional DOFs in five residues (GLU, GLN, ASN, ASP and ARG) were given significantly different definition of torsional states as shown in Fig. S3. The point here is not to raise an argument regarding the optimal way of assigning torsional states for side chain torsions. Rather, the differences give us an opportunity to test the sensitivity of the

ideal gas hypothesis on the specific definition of basic conformers. We repeated the analysis for RCONF2s. Corresponding N_{rconf2} and δS_{rconf2} , which are different from N_{rconf} and δS_{rconf} for some trajectory sets, were shown in Table 1. Nonetheless, the relationship between S_{rconf2}^{O-} and S_{rconf2}^{O-ig} , and that between F and S_{rconf2}^{O-ig} is essentially the same as what observed for RCONFs, the results are shown only for HEWL-b trajectory set (See Fig. 6.). These observations indicate that the ideal gas hypothesis is not sensitive to details in the specific definition of basic conformers.

Discussions

We did not include backbone torsional DOFs in the definition of RCONF(2)s based on the following consideration. Firstly, in protein folding, design and docking studies, backbone DOFs are usually treated explicitly, and free energy difference between two given backbone configurations are estimated with various scoring functions. Consistent with the idea that folded proteins have solid like backbone and liquid like side chains,³² it is found in our previous backbone conformational analysis of HEWL that the number of statistically significant combinations of backbone torsional states is very limited.²⁸ How to pick right backbone configurations out of astronomically large possible number of which for explicit free energy estimation is another difficult task to tackle in predictive tasks such as folding, design and docking, and we are actively investigating this issue. In our conformational partitions based on individual backbone dihedrals, backbone DOFs are not limited except for the one on which the projection is performed. No explicit restriction of backbone torsional DOFs is imposed on conformational partitions based on radius of gyration R_g and on number of native contacts N_{nc} . The observed validity of the ideal gas hypothesis indicates that change of backbone torsional states are effectively reflected by change of relevant side chain torsional states, at least in a statistical sense.

$S_{rconf-vib}$ in equation 14 includes contributions from both bonding and bending vibrational entropies, and from local part of torsional entropy within each RCONF. Seemingly, each RCONF allow significant local torsional motion since each defining torsion angle have a ~ 120 (or ~ 60

in a few cases) degree range to fluctuate. However, analysis of high resolution trajectories in sets HEWL-b and BPTI-c indicate this is not the case, the average life time of RCONFs are $\sim 2.3fs$ and $\sim 2.1fs$ respectively. This is due to the fact that even a small local torsional motion in one side chain torsional DOF is likely to be accompanied by change of torsional state in some other side chain torsions, the large number of side chain torsional DOFs results in short life time of RCONFs and correspondingly highly limited local torsional motion in each RCONF. In this study, we focused exclusively on conformational entropy based on RCONFs. The corresponding vibrational contributions to the configurational entropy ($S_{rconf-vib}$), together with other essential components of free energy in equations (2-3), were not analyzed. It is certainly desirable to have the capability to nail down these terms with high level of confidence. Unfortunately, reliable calculation of $S_{rconf-vib}$ is difficult by quasiharmonic or correlation based methods for most conformers due to limited number of snapshots available. No effective methods is presently available for calculation of correlation entropies between different type of molecular components (S_{pl} , S_{pv} and S_{lv}) in protein-ligand and other similar type of systems.

At first sight, the effective equivalence between change of conformational entropy and change of free energy seems exotic, and one would wondering what happened to enthalpic contributions and vibrational entropic contributions. It is important to note that the seemingly only important quantity N_{rconf} is a function of underlying molecular interactions including both intramolecular interactions within a protein molecule and intermolecular interactions between protein and solvent. Additionally, we emphasize that the constant local integral assumption does not necessarily limit $S_{rconf-vib}$, which might vary significantly for different RCONFs. We speculate that the silence of $S_{rconf-vib}$ in the observed change of free energy might due to its correlation, and consequently canceling effects, with other complex terms in equations 2 and 3.

The proposed end point free energy estimation methodology, as shown in equation 23, is in principle complementary to presently widely utilized methods such as Linear Interaction Energy (LIE) model³³ and MM/P(G)BSA,³⁴ especially for the cases where virtually no overlapping of conformational space exist for two end macrostates. Another advantage of equation 23 is that there

is no system dependent parameters to construct. It was demonstrated that side chain conformational entropies for given backbone configurations are not sensitive to force fields details.²⁹ However, the reported results are restricted to backbone of folded proteins or decoys that have reasonable packing density and surface exposure. To use equation 23 alone for selecting proper backbone configurations, the quality of solvation model is likely to be of critical importance.

United atom model is widely utilized to improve computational efficiency.^{35–37} Since RCONFs (or RCONF2s) are defined by heavy atoms, properly parameterized united atom models may potentially be utilized for counting number of RCONFs without significantly compromising accuracy. The present study is limited to native globular proteins. The physiological importance of membrane proteins and inherently disorder proteins are well acknowledged, some proteins interconvert between folded and unfolded states many times during their physiological life time. We are working on the generalization of the ideal gas hypothesis to these widely different scenarios, and to other complex molecular systems as well.

Conclusion

In summary, we proposed the ideal gas hypothesis to deal with lack of fundamental microstates in defining classical entropy of proteins. By utilizing the expediency of extensive MD trajectories in analyzing arbitrary partitions of protein conformational space, we tested the ideal gas hypothesis of RCONFs for a few globular proteins in native ensemble. The ideal gas hypothesis, while performs poorly for estimating absolute value of conformational entropy when the time resolution of observation is sufficiently fine, is demonstrated to be consistently robust as far as change of conformational entropy (or free energy) is concerned. A new end point free energy estimation method, which is a direct result of the ideal gas hypothesis, is also examined. This alternative free energy estimation scheme is applicable to cases where end states do not overlap in conformational space, which are highly challenging situations for presently available free energy methodologies.

Methods

BPTI-a trajectories were provided by DE Shaw.³¹ MD trajectories of HEWL (collectively 200 μ s comprising 2000 100-ns trajectories) were taken from our previous simulation study.²⁸ BPTI-b and KLKA trajectories (with pico-second resolution) were taken from another previous study.³⁸ BamC, BamE and CDK2 trajectories were generated in our group and details of these trajectories will be published in the future. For BamC, structure with PDB code 3TGO was solvated in 13736 water molecules, 39 Cl^- and 43 Na^+ ions. For BamE, structure with PDB code 2YH9 was solvated with 7391 water molecules, 22 Cl^- and 21 Na^+ ions. CHARMM22 force fields are used for simulations of BamC and BamE. CDK2 trajectories are based on AMBER ff12 force fields, 66 crystal structures (1FIN, 1GZ8, 1HCK, 1JST, 1JST, 1PF8, 1PW2, 1PXI, 1PXJ, 1PKX, 1PKM, 1W8C, 1Y8Y, 2A4L, 2B54, 2BPM, 2BPM, 2C4G, 2C4G, 2C5N, 2C5N, 2C5O, 2C5O, 2C5V, 2C5X, 2C69, 2C6K, 2C6L, 2CLX, 2EXM, 2UUE, 2V22, 2VTL, 2VTM, 2VTR, 2WEV, 2WFY, 2WHB, 2WIH, 2WIH, 2WPA, 2WXV, 2WXV, 2X1N, 2X1N, 3F5X, 3F5X, 3IGG, 3LFQ, 3PXF, 3PXZ, 3QHR, 3QQF, 3QQJ, 3RK9, 3RKB, 3S0O, 3UNK, 3WBL, 4BCK, 4EZ7, 4GCJ, 4I3Z, 4II5, 4KD1) are utilized to start 66 trajectories after each was solvated with 13851 water molecules, 49 Cl^- and 45 Na^+ ions, and production runs of \sim 200-ns are performed for each CDK2 system after equilibration. The same equilibration procedures as used in a previous study³⁸ was utilized for equilibration of these three protein simulation systems. The starting structural state of all 10-fs resolution trajectories are uniformly picked from corresponding pico-second(s) resolution trajectories. Specifically, 21101, 6124 and 4204 10-ps trajectories are generated for sets HEWL-b, BPTI-c and KLKA-b.

Acknowledgement

This research was supported by National Natural Science Foundation of China under grant number 31270758. Computational resources were partially supported by High Performance Computing Center of Jilin University, China. We thank DE Shaw Research for providing BPTI trajectories.

We thank Zhonghan Hu for critical reading of the manuscript.

References

- (1) Jaynes, E. T. *Phys. Rev.* **1957**, *106*, 620–630.
- (2) Jaynes, E. T. *Phys. Rev.* **1957**, *108*, 171–190.
- (3) Wehrl, A. *Reviews of Modern Physics* **1978**, *50*, 221–260.
- (4) Frederick, K. K.; Marlow, M. S.; Valentine, K. G.; Wand, A. J. *Nature* **2007**, *448*, 325–329.
- (5) Tzeng, S.-R.; Kalodimos, C. G. *Nature* **2012**, *19*.
- (6) Kasinath, V.; Sharp, K. a.; Wand, a. J. *Journal of the American Chemical Society* **2013**, *135*, 15092–100.
- (7) Wand, A. J. *Current Opinion in Structural Biology* **2013**, *23*, 75 – 81, Folding and binding / Protein-nucleic acid interactions.
- (8) Reinhard, F.; Grubmüller, H. *The Journal of chemical physics* **2007**, *126*, 014102.
- (9) Tyka, M. D.; Sessions, R. B.; Clarke, A. R. *The journal of physical chemistry. B* **2007**, *111*, 9571–80.
- (10) Reinhard, F.; Lange, O. F.; Hub, J. S.; Haas, J.; Grubmüller, H. *Computer Physics Communications* **2009**, *180*, 455–458.
- (11) Wang, L.; Abel, R.; Friesner, R. a.; Berne, B. J. *Journal of chemical theory and computation* **2009**, *5*, 1462–1473.
- (12) Gerogiokas, G.; Calabro, G.; Henchman, R. H.; Southey, M. W. Y.; Law, R. J.; Michel, J. *Journal of Chemical Theory and Computation* **2014**, *10*, 35–48.

- (13) Suárez, D.; Díaz, N. *Wiley Interdisciplinary Reviews: Computational Molecular Science* **2014**, 5, 1–26.
- (14) Falconer, R. J.; Penkova, A.; Jelesarov, I.; Collins, B. M. *Journal of molecular recognition : JMR* **2010**, 23, 395–413.
- (15) Lipari, G.; Szabo, A. *Journal of the American Chemical Society* **1982**, 104, 4546–4559.
- (16) Lipari, G.; Szabo, A. *Journal of the American Chemical Society* **1982**, 104, 4559–4570.
- (17) Karplus, M.; Ichiye, T.; Pertirr, B. M. *Biophysical Journal* **1987**, 52, 1083–1085.
- (18) Chang, C.-e.; Gilson, M. K. *Journal of American Chemical Society* **2004**, 126, 13156–13164.
- (19) Chang, C.-e. a.; Chen, W.; Gilson, M. K. *Proceedings of the National Academy of Sciences of the United States of America* **2007**, 104, 1534–9.
- (20) Numata, J.; Knapp, E.-W. *Journal of Chemical Theory and Computation* **2012**, 8, 1235–1245.
- (21) Bower, M. J.; Cohen, F. E.; Dunbrack, R. L. *Journal of molecular biology* **1997**, 267, 1268–1282.
- (22) Shapovalov, M. V.; Dunbrack, R. L. *Structure* **2011**, 19, 844–858.
- (23) Scouras, A. D.; Daggett, V. *Protein Science* **2011**, 20, 341–352.
- (24) Li, D.-W.; Brüschweiler, R. *Physical Review Letters* **2009**, 102, 118108.
- (25) Shannon, C. *The Bell System Technical Journal* **1948**, 27, 379–423.
- (26) Cormen, T. H.; Leiserson, C. E.; Rivest, R. L.; Stein., C. *Introduction to Algorithms*, 3rd ed.; MIT Press and McGraw-Hill, 2009.
- (27) Csermely, P.; Palotai, R.; Nussinov, R. *Trends in Biochemical Sciences* **2010**, 35, 539 – 546.
- (28) Wang, K.; Long, S.; Tian, P. *submitted*

- (29) Zhang, J.; Liu, J. S. *PLoS computational biology* **2006**, *2*, e168.
- (30) Zhang, J.; Chen, R.; Tang, C.; Liang, J. *The Journal of Chemical Physics* **2003**, *118*, 6102.
- (31) Shaw, D. E.; Maragakis, P.; Lindorff-Larsen, K.; Piana, S.; Dror, R. O.; Eastwood, M. P.; Bank, J. a.; Jumper, J. M.; Salmon, J. K.; Shan, Y.; Wriggers, W. *Science (New York, N.Y.)* **2010**, *330*, 341–6.
- (32) Lindorff-Larsen, K.; Best, R. B.; DePristo, M. A.; Dobson, C. M.; Vendruscolo, M. *Nature* **2005**, *433*, 128–132.
- (33) Hansson, T.; Marelius, J.; Åřqvist, J. *Journal of Computer-Aided Molecular Design* **1998**, *12*, 27–35.
- (34) Kuhn, B.; Kollman, P. A. *Journal of Medicinal Chemistry* **2000**, *43*, 3786–3791.
- (35) Berger, O.; Edholm, O.; Jähnig, F. *Biophysical Journal* **1997**, *72*, 2002 – 2013.
- (36) Chiu, S.-W.; Pandit, S. A.; Scott, H. L.; Jakobsson, E. *The Journal of Physical Chemistry B* **2009**, *113*, 2748–2763, PMID: 19708111.
- (37) Tjörnhammar, R.; Edholm, O. *Journal of Chemical Theory and Computation* **2014**, *10*, 5706–5715.
- (38) LI Wenzhao, T. P., MENG Wei *CHEMICAL RESEARCH IN CHINESE UNIVERSITIES* **2015**, *31*, 149.

Table 1: The list of the studied protein MD trajectory sets. For the trajectory sets of the same protein with different resolutions, we use *-a*, *-b* and *-c* to make the distinction. N_{res} : number of residues; N_{tor} : number of heavy-atom side chain torsional DOFs utilized in defining RCONF(2)s; Δt : time interval for saving MD snapshots in *ps*; N_{snap} : total number of snapshots in the given trajectory set; N_{rconf} : total number of RCONFs in the given trajectory set; N_{rconf2} : total number of RCONF2s in the given trajectory set; δS_{rconf} (in the unit of k_B): deviation of RCONF entropy from the ideal gas value as indicated in equation 34; δS_{rconf2} (in the unit of k_B): deviation of RCONF2 entropy from the ideal gas value, defined similarly with δS_{rconf} .

Protein	N_{res}	Δt	N_{snap}	N_{tor}	N_{rconf}	δS_{rconf}	N_{rconf2}	δS_{rconf2}
BPTI-a	58	250	4120838	113	4120838	0.00	4120815	0.000007
HEWL	129	4	50000000	196	49990573	0.000262	49668320	0.009333
CDK2	298	2	6560590	565	6560590	0.00	6560590	0.00
BamC	190	2	943505	339	943505	0.00	921630	0.033388
BamE	68	2	1679968	120	1674297	0.004708	1094439	0.639567
KLKA	223	1	3741963	413	3741963	0.009684	3653756	0.034387
BPTI-b	58	1	3560127	113	3523927	0.014257	1425354	1.330897
KLKA-b	223	0.01	4204000	413	4201317	0.000884	823555	2.326272
HEWL-b	129	0.01	21101000	196	9194924	1.214137	2845936	2.837891
BPTI-c	58	0.01	6124000	113	2923546	1.074068	76140	5.477714

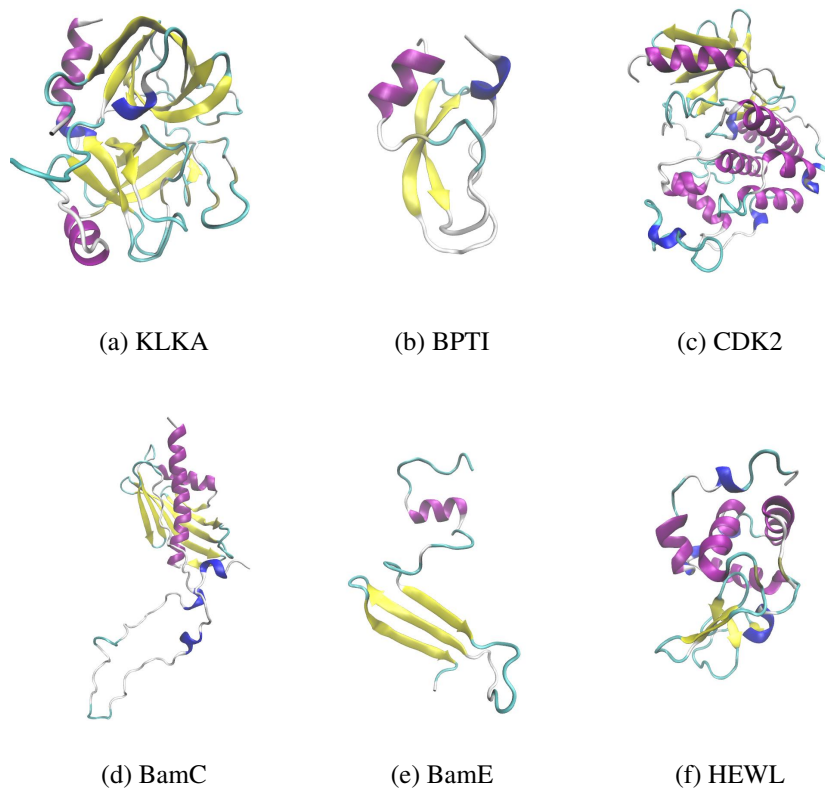


Figure 1: Structures of proteins utilized in this study. α helices are in purple, β strands are in yellow, 3-10 helices are in blue, loops are in cyan and turns are in white. KLKA (Porcine Pancreatic Kallikrein A); BPTI (Bovine Pancreatic Trypsin Inhibitor); CDK2 (Cyclin-dependent kinase 2); Bam (The β -barrel assembly machine); HEWL (Hen Egg White Lysozyme). Graphics are prepared using VMD.

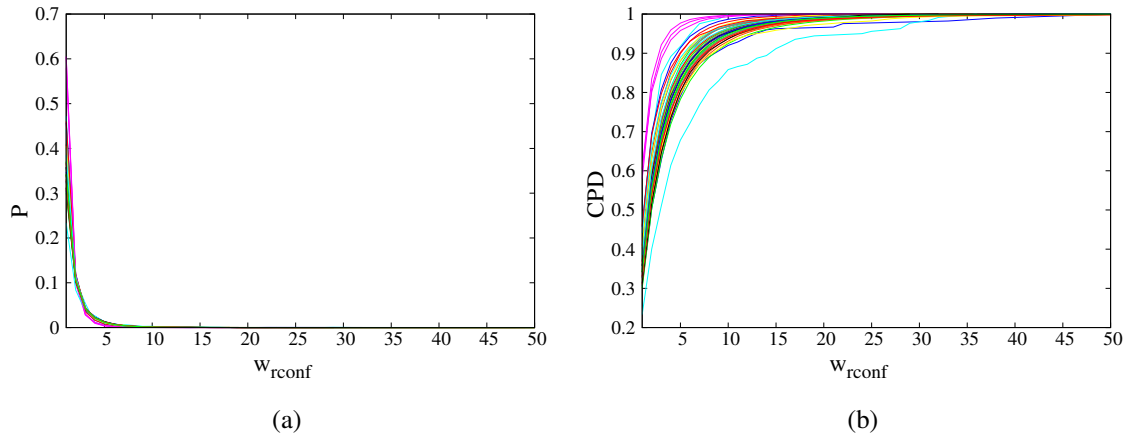


Figure 2: Statistical weight of RCONFs (w_{rconf}) as represented by the number of snapshots in 10-fs interval HEWL trajectories. a) Probability (P) of various w_{rconf} in 37 representative conformational partitions. b) Cumulative probability density (CPD) of w_{rconf} for the same set of conformational partitions as in a).

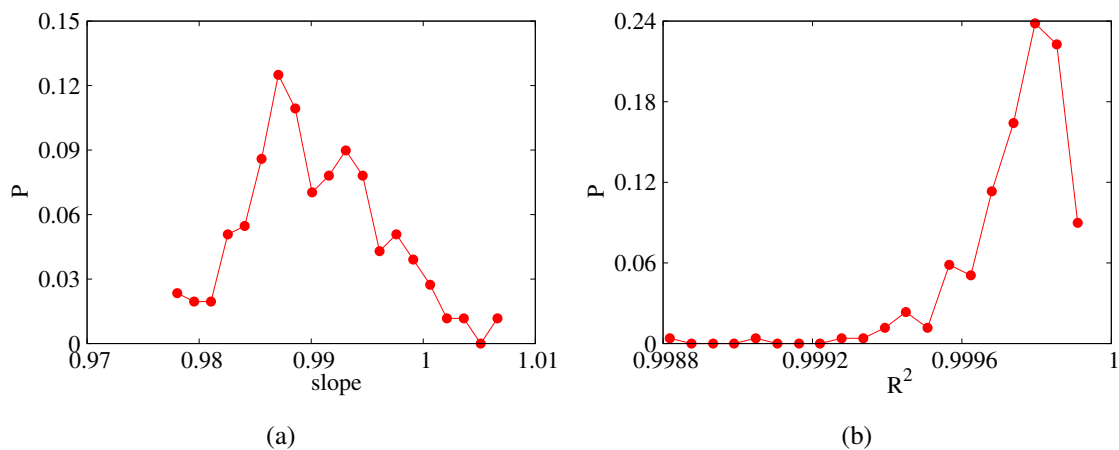


Figure 3: Probability (P) distributions of a)slopes and b)squared linear correlation coefficients calculated from 256 S_{rconf}^{O-} vs. S_{rconf}^{O-ig} plots based on 256 different ways of conformational partitions performed on HEWL-b trajectory set. Each way of conformational partition corresponds to projection of all snapshots onto a given backbone dihedral (ϕ or ψ).

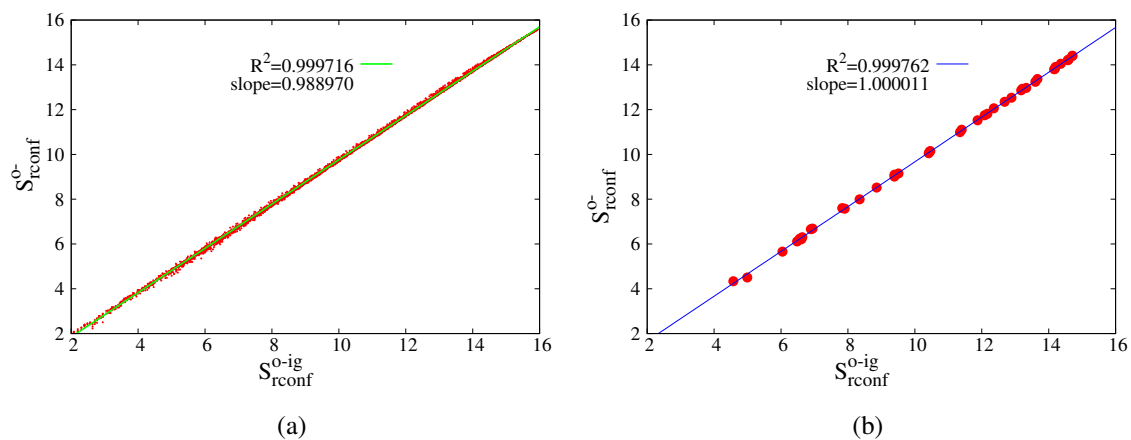


Figure 4: S_{rconf}^{O-} vs. S_{rconf}^{O-ig} (both in the unit of k_B) plots of HEWL-b for two different types of overlapping conformational partitions. a) For 5120 partitions generated from projection onto 256 backbone dihedrals. b) For 40 partitions generated from projection onto radius of gyration R_g and the number of native contacts N_{nc} .

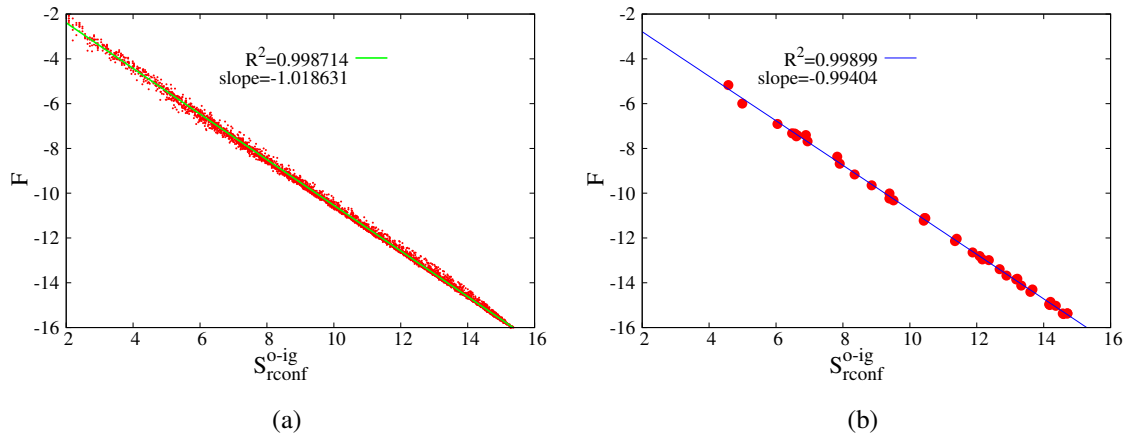


Figure 5: Relative free energy F (in the unit of $k_B T$) vs. S_{rconf}^{O-ig} (in the unit of k_B) plots for two different types of overlapping conformational partitions. a) For 5120 partitions generated from projection onto 256 backbone dihedrals. b) For 40 partitions generated from projection onto radius of gyration R_g and the number of native contacts N_{nc} .

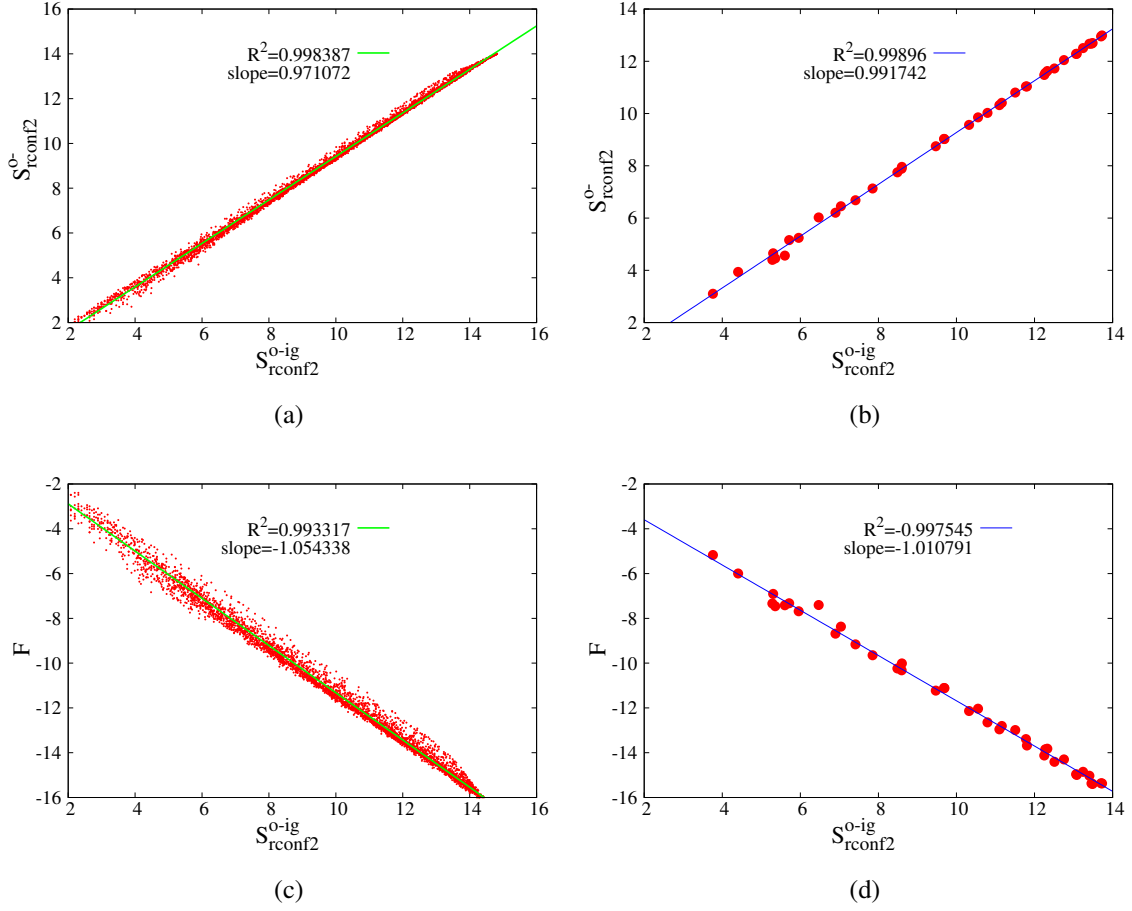


Figure 6: Results for RCONF2 based analysis of HEWL-b trajectory set. S_{rconf}^{O-} vs. S_{rconf}^{O-ig} (both in the unit of k_B) plots for a) 5120 partitions generated from projection onto 256 backbone dihedrals and b) 40 partitions generated from projection onto radius of gyration R_g and the number of native contacts N_{nc} . c) and d) Relative free energy F (in the unit of $k_B T$) vs. S_{rconf}^{O-ig} (in the unit of k_B) plots for the same two sets of the conformational partitions.